

Q1

e) Explain in a few sentences what happens when the model's `train_op` is called (what gets computed during forward propagation, what gets computed during backpropagation, and what will have changed after the op has been run?).

The loss and the gradients of all `tf.Variable()`'s are computed in forward prop and back prop , respectively. The `tf.Variable()`'s are updated after each iteration.

Q2

a)

Iteration	Stack	Buffer	New Dependency	Transition
0	ROOT	I parsed this sentence correctly		
1	ROOT I	Parsed this sentence correctly		Shift
2	ROOT I parsed	This sentence correctly		Shift
3	ROOT parsed	This sentence correctly	Parsed -> I	Left-Arc
4	ROOT parsed this	Sentence correctly		Shift
5	ROOT parsed this sentence	correctly		Shift
6	ROOT parsed sentence	Correctly	Sentence -> this	Left-arc
7	ROOT parsed	Correctly	Parsed -> sentence	Right-arc
8	ROOT parsed correctly			Shift
9	ROOT parsed		Parsed -> correctly	Right-arc
10	ROOT		ROOT -> parsed	Right-arc

b) A sentence containing n words will be parsed in how many steps (in terms of n)? Briefly explain why.

A sentence with n words will be parsed in $2n$ steps, since each word needs a SHIFT and an ARC step.

d) When implementing the mini-batch parser, shallow copies of the Parser list should be created. That way, when the self.dependencies of the unfinished Parser list are updated, the same is done for the original Parser list. Also, removing completed Parser objects in the unfinished Parser list will not affect the original this way.

Also, to remove completed Parsers in the unfinished Parser list, first create a shallow copy of a list of the first `batch_size` Parser objects, so that removing the completed Parser objects will not trigger an `IndexError`.

Use `parse_step` instead of `parse`, since `parse` performs multiple parsing with a list of transitions. If you pass in only one transition, it treats that transition as a list which means each `parse_step` in `parse` takes in a character from the transition. If you used an `else` for the last transition command (S, LA, RA), the parser will perform the last transition for each character that it misinterprets as a transition.

f)

(2 points, written) We will regularize our network by applying Dropout³. During training this randomly sets units in the hidden layer \mathbf{h} to zero with probability p_{drop} and then multiplies \mathbf{h} by a constant γ (dropping different units each minibatch). We can write this as

$$\mathbf{h}_{drop} = \gamma \mathbf{d} \circ \mathbf{h}$$

where $\mathbf{d} \in \{0, 1\}^{D_h}$ (D_h is the size of \mathbf{h}) is a mask vector where each entry is 0 with probability p_{drop} and 1 with probability $(1 - p_{drop})$. γ is chosen such that the value of \mathbf{h}_{drop} in expectation equals \mathbf{h} :

$$\mathbb{E}_{p_{drop}}[\mathbf{h}_{drop}]_i = h_i$$

for all $0 < i < D_h$. What must γ equal in terms of p_{drop} ? Briefly justify your answer.

$$\begin{aligned} E_{p_{drop}}[h_{drop}] &= \sum p_{drop} h_{drop} \\ &= p(h_{drop} = 0)(h_{drop} = 0) + p(h_{drop} = \gamma h)(h_{drop} = \gamma h) \\ &= p_{drop}(0) + (1 - p_{drop})(\gamma h) \\ \therefore \gamma &= \frac{1}{1 - p_{drop}} \end{aligned}$$

g)

(i) First, Adam uses a trick called *momentum* by keeping track of \mathbf{m} , a rolling average of the gradients:

$$\begin{aligned} \mathbf{m} &\leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\theta} J_{minibatch}(\theta) \\ \theta &\leftarrow \theta - \alpha \mathbf{m} \end{aligned}$$

where β_1 is a hyperparameter between 0 and 1 (often set to 0.9). Briefly explain (you don't need to prove mathematically, just give an intuition) how using \mathbf{m} stops the updates from varying as much. Why might this help with learning?

Momentum prevents the gradients from directly influencing/changing the parameters, since it is a moving average of the gradients which smoothens the noise and the gradients of the current iteration are reduced by a factor of 10. Hence there is less fluctuation in the updates due to less noise from the gradients and the current gradient has less magnitude which makes it less influential.

It helps learning because the gradients at a certain iteration might be 0, which causes naïve optimizers like SGD to stop at that local minimum even if the loss at that point is high. With momentum, the update to the parameters are not zero even if the gradient is zero, which prevents the parameters from being stuck at a bad local minimum by having that residual momentum shift the parameters away from the local minimum.

- (ii) Adam also uses *adaptive learning rates* by keeping track of \mathbf{v} , a rolling average of the magnitudes of the gradients:

$$\begin{aligned}\mathbf{m} &\leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta}) \\ \mathbf{v} &\leftarrow \beta_2 \mathbf{v} + (1 - \beta_2) (\nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta}) \circ \nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta})) \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \alpha \circ \mathbf{m} / \sqrt{\mathbf{v}}\end{aligned}$$

where \circ and $/$ denote elementwise multiplication and division (so $\mathbf{z} \circ \mathbf{z}$ is elementwise squaring) and β_2 is a hyperparameter between 0 and 1 (often set to 0.99). Since Adam divides the update by $\sqrt{\mathbf{v}}$, which of the model parameters will get larger updates? Why might this help with learning?

Parameters with small gradients will have larger updates, since small gradients mean that they are converging to a possibly bad local minimum. By increasing the update magnitude, it can help these parameters to shift away from the local minimum.

h)

Baseline: Best UAS on dev: 87.94, UAS on test: 88.48

With L2-regularization: Lower

With additional hidden layer of same dimensions: Lower

Q3

- (a) (4 points, written) Conventionally, when reporting performance of a language model, we evaluate on *perplexity*, which is defined as:

$$PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{\bar{P}(\mathbf{x}_{\text{pred}}^{(t+1)} = \mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}}$$

i.e. the inverse probability of the correct word, according to the model distribution \bar{P} .

- (i) (2 points) Show how you can derive perplexity from the cross-entropy loss (*Hint: remember that $\mathbf{y}^{(t)}$ is one-hot!*). **This should be a very short problem - not too perplexing!**
- (ii) (1 point) Now use this relationship between perplexity and cross-entropy to show that minimizing the geometric mean perplexity, $\left(\prod_{t=1}^T PP(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})\right)^{1/T}$, is equivalent to minimizing the arithmetic mean cross-entropy loss, $\frac{1}{T} \sum_{t=1}^T CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})$, across the training set. (*Hint: for any positive function f , minimizing $\log(f)$ is equivalent to minimizing f itself*)
- (iii) (1 point) Suppose you have a vocabulary of $|V|$ words and your “language model” works by picking the next word uniformly at random from the vocabulary (mathematically, $\bar{P}(\mathbf{x}^{(t+1)} = \mathbf{w}_j \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = 1/|V|$ for every word \mathbf{w}_j in the vocabulary). What would the perplexity $PP(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})$ for a single word to be? Compute the corresponding cross-entropy loss when $|V| = 10000$.

$$\begin{aligned}\text{(i)} \quad CE(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)}) &= -\sum_{j=1}^{|V|} y_j^{(t)} \log(\hat{y}_j^{(t)}) \\ e^{CE(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)})} &= e^{-\sum_{j=1}^{|V|} y_j^{(t)} \log(\hat{y}_j^{(t)})} = \frac{1}{e^{\sum_{j=1}^{|V|} y_j^{(t)} \log(\hat{y}_j^{(t)})}} = \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} e^{\log(\hat{y}_j^{(t)})}} \\ &= \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \hat{y}_j^{(t)}} = PP(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)})\end{aligned}$$

$$(ii) \quad \min_x \log(f(x)) = \min_x f(x)$$

$$\begin{aligned} \min_y \left[\prod_{t=1}^T PP(\hat{y}^{(t)}, y^{(t)}) \right]^{1/T} &= \min_y \log \left[\prod_{t=1}^T PP(\hat{y}^{(t)}, y^{(t)}) \right]^{1/T} \\ &= \min_y \frac{1}{T} \log \left(\prod_{t=1}^T PP(\hat{y}^{(t)}, y^{(t)}) \right) \\ &= \min_y \frac{1}{T} \sum_{t=1}^T \log(PP(\hat{y}^{(t)}, y^{(t)})) = \min_y \frac{1}{T} \sum_{t=1}^T CE(\hat{y}^{(t)}, y^{(t)}) \end{aligned}$$

$$(iii) \quad \text{Let } \hat{y}_j^{(t)} = \frac{1}{|V|} \quad \forall j, \quad PP(\hat{y}^{(t)}, y^{(t)}) = \frac{1}{1(\frac{1}{|V|})+0} = |V|$$

$$\text{If } |V| = 10000, \quad CE(\hat{y}^{(t)}, y^{(t)}) = \log(PP(\hat{y}^{(t)}, y^{(t)})) = \log(10000) = 4$$

- (b) (7 points, written) Compute the gradients of the loss J (the cross-entropy loss) with respect to the following model parameters at a single point in time t (to save a bit of time, you don't have to compute the gradients with the respect to the biases \mathbf{b}_1 and \mathbf{b}_2):

$$\frac{\partial J^{(t)}}{\partial U} \quad \frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t)}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big|_{(t)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(t)}$$

where $\Big|_{(t)}$ denotes the gradient for the appearance of that parameter at time t (in other words, $\mathbf{h}^{(t-1)}$ is taken to be fixed, so you don't need to take it's derivative when applying the chain rule and backpropagate to earlier timesteps - you'll do that in part (c)).

Additionally, compute the derivative with respect to the *previous* hidden layer value:

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}}$$

Hint: you may want to define the following variables:

$$\begin{aligned} \mathbf{z}^{(t)} &= \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1 \\ \boldsymbol{\theta}^{(t)} &= \mathbf{U} \mathbf{h}^{(t)} + \mathbf{b}_2 \end{aligned}$$

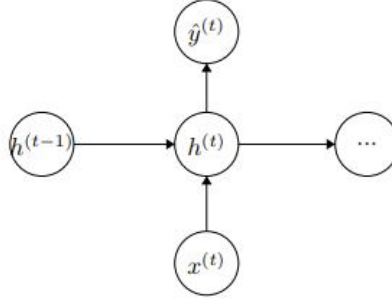
and define (and compute the value of) the following error terms:

$$\begin{aligned} \delta_1^{(t)} &= \frac{\partial J}{\partial \boldsymbol{\theta}^{(t)}} \\ \delta_2^{(t)} &= \frac{\partial J}{\partial \mathbf{z}^{(t)}} = \delta_1^{(t)} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} \end{aligned}$$

$$\begin{aligned} \frac{\partial J^{(t)}}{\partial U} &= \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial U} = \delta_1^{(t)} [\mathbf{h}^{(t)}]^T & \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big|_t &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_e} \Big|_t = \delta_2^{(t)} [\mathbf{e}^{(t)}]^T \\ \frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t)}} &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{e}^{(t)}} = \mathbf{W}_e^T \delta_2^{(t)} & \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_t &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_h} \Big|_t = \delta_2^{(t)} [\mathbf{h}^{(t-1)}]^T \\ \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \frac{\partial J^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \mathbf{W}_h^T \delta_2^{(t)} \end{aligned}$$

$$\text{Where} \quad \delta_1^{(t)} = \hat{y}^{(t)} - y^{(t)} \quad \delta_2^{(t)} = U^T \delta_1^{(t)} \circ [\sigma(\mathbf{z}^{(t)}) \circ [1 - \sigma(\mathbf{z}^{(t)})]]$$

(c) (7 points, written) Below is a sketch of the network at a single timestep:

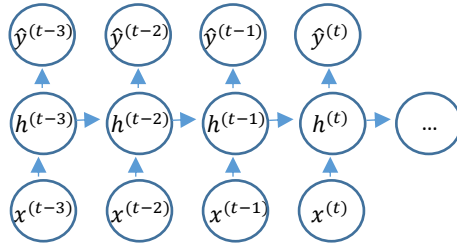


Draw the “unrolled” network for 3 timesteps (your picture should show $h^{(t-3)}$ to $h^{(t)}$). Next, compute the following backpropagation-through-time gradients:

$$\frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t-1)}} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_e} \Big|_{(t-1)} \quad \frac{\partial J^{(t)}}{\partial \mathbf{W}_h} \Big|_{(t-1)}$$

where $\Big|_{(t-1)}$ denotes the gradient for the appearance of that parameter at time $(t-1)$. Use the error term $\gamma^{(t-1)} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}}$ to represent the derivative you computed in the previous part (it should show up in three gradients you compute).

We have to compute multiple gradients with respect to parameters like \mathbf{W}_e because they are used multiple times in the feed-forward computation. The final gradient for the parameters will be the sum of all their gradients over time (e.g., the gradient $\frac{\partial J^{(t)}}{\partial \mathbf{W}_e} = \sum_{i=0}^t \frac{\partial J^{(i)}}{\partial \mathbf{W}_e} \Big|_{(i)}$).



Let $\gamma^{(t-1)} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \mathbf{W}_h^T \delta_2^{(t)}$

$$\frac{\partial J}{\partial \mathbf{e}^{(t-1)}} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{e}^{(t-1)}} = \mathbf{W}_e^T (\gamma^{(t-1)} \circ \sigma'(z^{(t-1)}))$$

$$\frac{\partial J}{\partial \mathbf{W}_e} \Big|_{t-1} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{W}_e} = (\gamma^{(t-1)} \circ \sigma'(z^{(t-1)})) [\mathbf{e}^{(t-1)}]^T$$

$$\frac{\partial J}{\partial \mathbf{W}_h} \Big|_{t-1} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{W}_h} = (\gamma^{(t-1)} \circ \sigma'(z^{(t-1)})) [\mathbf{h}^{(t-2)}]^T$$

- (d) (3 points, written) Given $\mathbf{h}^{(t-1)}$, how many operations are required to perform backpropagation for a single timestep (i.e., compute the terms you found in part (b)). Express your answer in big-O notation in terms of the dimensions d , D_h and $|V|$ (Equation 1). Don't worry about the gradients you didn't compute (with respect to \mathbf{b}_1 and \mathbf{b}_2), they actually don't change the answer!

(Hint: You only have to worry about matrix multiplications, the other operations do not change the big-O runtime. Multiplying a vector by an n by m matrix takes $O(nm)$ operations.)

$$\frac{\partial J^{(t)}}{\partial U} = \delta_1^{(t)} [h^{(t)}]^T, \delta_1^{(t)} = \hat{y}^{(t)} - y^{(t)} \rightarrow O(|V|), [h^{(t)}]^T \rightarrow O(D_h)$$

$$\therefore \frac{\partial J^{(t)}}{\partial U} \rightarrow O(|V| \cdot D_h)$$

$$\delta_2^{(t)} = U^T \delta_1^{(t)} \circ [\sigma(z^{(t)}) \circ [1 - \sigma(z^{(t)})]] = U^T \delta_1^{(t)} \circ \sigma'(z^{(t)})$$

$$U^T \delta_1^{(t)} \rightarrow O(|V| \cdot D_h \cdot 1), \sigma(z^{(t)}) \rightarrow 3O(D_h) \rightarrow O(D_h), \sigma'(z^{(t)}) \rightarrow O(D_h^2)$$

$$[U^T \delta_1^{(t)}] \circ \sigma'(z^{(t)}) \rightarrow O(D_h^2)$$

$$\therefore \delta_2^{(t)} \rightarrow O(|V| \cdot D_h), \text{ Assuming } |V| > D_h$$

$$\frac{\partial J^{(t)}}{\partial e^{(t)}} = W_e^T \delta_2^{(t)} \rightarrow O(d \cdot D_h) \therefore \frac{\partial J^{(t)}}{\partial e^{(t)}} \rightarrow O(|V| \cdot D_h), \text{ Assuming } |V| > d$$

$$\frac{\partial J^{(t)}}{\partial h^{(t-1)}} = W_h^T \delta_2^{(t)} \rightarrow O(D_h^2) \therefore \frac{\partial J^{(t)}}{\partial h^{(t-1)}} \rightarrow O(|V| \cdot D_h)$$

$$\left. \frac{\partial J^{(t)}}{\partial W_e} \right|_t = \delta_2^{(t)} [e^{(t)}]^T \rightarrow O(d \cdot D_h) \therefore \left. \frac{\partial J^{(t)}}{\partial W_e} \right|_t \rightarrow O(|V| \cdot D_h)$$

$$\left. \frac{\partial J^{(t)}}{\partial W_h} \right|_t = \delta_2^{(t)} [h^{(t-1)}]^T \rightarrow O(D_h^2) \therefore \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_t \rightarrow O(|V| \cdot D_h)$$

The biases b_1, b_2 do not change the complexity since they are $O(D_h), O(|V|)$ respectively. We drop lower order complexities.

Solution: $O(|V| \cdot D_h + d \cdot D_h + D_h^2)$ for the general case where the values of $|V|, D_h$ and d are unknown.

- (e) (3 points, written) Now suppose you have a sequence of T words. How many operations are required to compute the gradient of the loss with respect to the model parameters across the entire sequence $(\sum_{t=1}^T J^{(t)}(\theta))$? Assume we backpropagate through time all the way to $t = 0$ for each word.

Hint: Look at the computational graph you drew in part (c). Will we have to pass an error signal (upstream gradient) into $h^{(0)}$ T times (once for the loss from every word), or can we be more efficient than that? Therefore, how many times will we have to do the single timestep (your answer to (d)) computations?

Computing $\frac{\partial J^{(t)}}{\partial *}$ requires T steps for each t , so it is $O(T[|V| \cdot D_h + d \cdot D_h + D_h^2])$. Only one backward pass is made since we sum the gradients $\frac{\partial J^{(t)}}{\partial *}$ up, else it would be $O(T^2[...])$. T is not a constant so we cannot drop it.

- (f) (1 point, written) Which term in your big-O expressions is likely the largest? Which layer in the RNN is that term from?

$O(|V| \cdot D_h)$ is largest. Due to $\delta_2^{(t)}$.